

Review

Summative OSCEs in undergraduate medical education

Gerry Gormley

Accepted 3 August 2011

INTRODUCTION

Making judgements on the competency of our peers and trainees is important in patient healthcare.¹ Inaccuracies in such judgements could place patients at risk. First described in 1979², Objective Structured Clinical Examinations (OSCEs) have become one of the most widely used methods of assessing aspects of clinical competency in healthcare education.³ This method of assessment was originally developed in order to address the unreliability and lack of generalisability of traditional forms of clinical assessment such as the *long case*.⁴ The overarching philosophy in OSCEs is that all candidates are presented with the same clinical tasks, to be completed in the same timeframe and are scored using structured marking schemes.² Compared to the *long case*, OSCEs reduce bias relating to the type of clinical case selected and who performs the assessment. Ideally the only variance in an OSCE should be the candidate's performance. In *formative* forms of assessment the main purpose is to provide feedback to the student. *Summative* forms of assessment define those who have achieved a passing standard and can progress in their studies.⁵ This article aims to provide a review of summative OSCEs in undergraduate medical education.

ASSESSMENT OF CLINICAL COMPETENCY: WHERE DO OSCEs FIT INTO THE BIGGER PICTURE?

The assessment of clinical competence is of significant importance. The General Medical Council emphasises the importance of accurately assessing the competency of medical students.⁶ Such decisions help to protect patients by determining whether candidates can progress to higher levels of study or medical qualification.

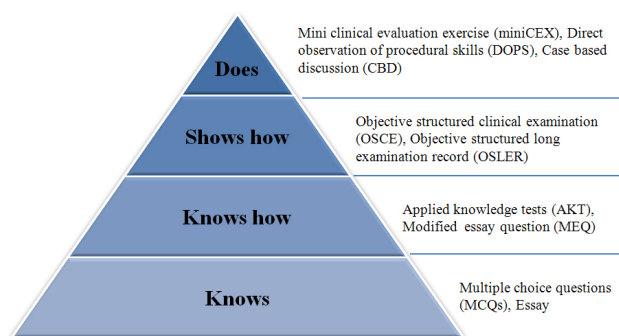


Fig 1. Adapted version of Millers' pyramid of clinical competency.

Miller provides a conceptual framework for assessing clinical competency (Figure 1).⁷ This pyramidal model describes the various domains of clinical competency. In achieving clinical competency, candidates are not only required to demonstrate that they *know* the facts which underpin clinical practice but also *know how* to apply these facts. Crucially they also need to *show* that they can perform the clinical tasks and skills. This facet of clinical competence relates more to behavioural than cognitive attributes. OSCEs are a common method of assessing the *shows how* aspects of clinical competency.

Despite the popularity of OSCEs it is important to note they do not provide a complete profile of an individual's level of competency. No valid single method of assessment exists. OSCEs aim to assess certain aspects of clinical competency. Using multiple assessment tools longitudinally is considered the best approach in forming a more holistic opinion on an individual's level of clinical competency.⁵ By using several methods of assessment the inadequacies of individual methods may be overcome.⁸ Attaining clinical competence is not a one-off event but a career long learning routine.⁵

WHAT IS THE TYPICAL FORMAT OF AN OSCE?

In the UK there is no standard operating procedure for running OSCEs. Therefore there will always be institutional variation in how OSCEs are delivered. However the underlying principles of OSCEs are common to all medical schools. In an OSCE, candidates sequentially rotate around a series of structured clinical cases or stations. Typically in a final year OSCE there may be anywhere between 10-20 individual stations. Stations aim to sample across a wide range of clinical competencies (Figure 2). For example:

- communication and professionalism skills (*e.g. breaking bad news*)
- history taking skills (*e.g. taking a history from a patient presenting with acute chest pain*)
- physical examination skills (*e.g. performing a respiratory examination*)
- clinical-reasoning skills (*e.g. interpreting clinical data and then prescribing therapy on a drug chart*)

Department of General Practice, Dunluce Health Centre, Belfast BT7 9HR.

Correspondence to Dr G Gormley

g.gormley@qub.ac.uk

- practical / technical skills (e.g. insertion of a peripheral venous cannula)

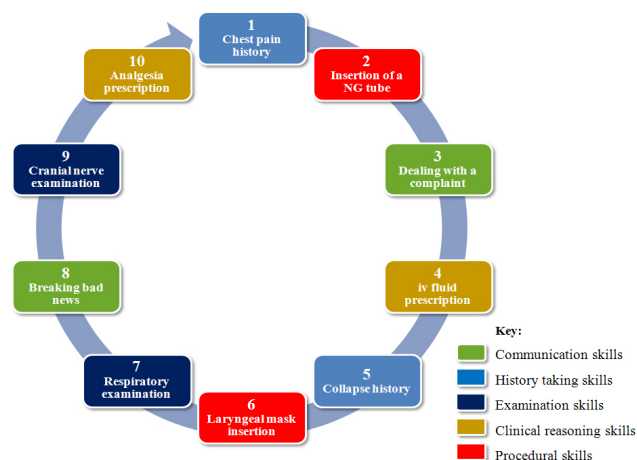


Fig 2. Graphical representation of a theoretical graduating medical OSCE.

At each station candidates are assigned a specific clinical task to perform. In these stations they may encounter a *real* or *simulated* patient, manikin, part-task manikin (i.e. a simulated patient in combination with a manikin), a computer based simulation (e.g. clinical video of a real patient with signs of Parkinson's disease) or clinical information (e.g. a fluid balance chart, blood results and an intravenous fluid prescription chart). Each station has a predefined structured marking scheme or *checklist*. There usually is an assessor in each station who observes the candidate and scores their performance according to the checklist. After a set time period, a bell will signal for candidates to move on to the next station. The circuit of stations is followed in sequence by all candidates. In circumstances where there are a large number of candidates, the OSCE may run across different examination venues and sometimes over the course of one day or more.

WHAT MAKES AN OSCE A GOOD FORM OF ASSESSMENT?

There are many attributes of a good and useful test. Van der Vleuten described five such criteria – namely: *reliability*, *validity*, *educational impact*, *cost efficiency* and *acceptability* of the test.⁹ Although excelling in all criteria would be ideal, pragmatically there often has to be compromise.

Reliability of OSCEs

Reliability of a test is a measure of its reproducibility and accuracy. In other words the degree to which a test consistently measures what it is intended to measure. OSCEs are widely considered to be a reliable form of assessment. There are many features of OSCEs that contribute to their reliability. Assessor consistency is improved by the use of highly structured marking schemes. Individual assessor bias is reduced by the use of multiple assessors. Ultimately having multiple cases, and sufficient test time, are the most important features that contribute to the reliability of OSCEs.¹⁰ Godfrey Pell and colleagues describe a number of metrics (*such as Cronbach's alpha and R² coefficient*) that give an indication of

the reliability and quality of an OSCE.¹¹ The GMC emphasise the importance of using such reliability metrics to quality assure and improve the assessment process.¹²

Validity of OSCEs

The validity of an OSCE is determined by its ability to actually measure what it is intended to measure. In other words an OSCE is considered valid if it succeeds in measuring competencies that it was originally designed to test. There are different types of validity evidence. For example *content validity* of an OSCE is a measure of how well the OSCE stations match the learning outcomes of the course. Blueprinting an OSCE (i.e. stations selected to be used in an OSCE are representatively and systematically sampled from the entire range of learning outcomes for the course) enhances its *content validity*.

Educational impact of OSCEs

Assessment provides a crucial role in the educational process. Not only does it check that learning has occurred but it can provide a powerful influence on future learning.⁸⁻¹⁰ The current emphasis in education is moving away from 'assessment of learning' to 'assessment for learning'. Strategically designing OSCE content and format can have both a positive and negative impact on students' learning behaviours.^{9,13}

Students often focus their studies on what they predict will occur in an OSCE. The challenge for faculty is to encourage students not to focus on predictions but the stated learning outcomes of the course. Such an effect is known as *consequential validity*. A criticism of OSCEs is that they can promote students to learn the *checklist* rather than having a deeper understanding of the skill.¹⁴ Given these concerns there is now a trend in more senior level OSCEs to group together single 'lower-level' *checklist* items to more 'higher-level' items – also known as "*chunking*".¹¹ For example instead of using separate single marks for hand washing, identification of patient, explaining purpose of encounter – these items are grouped into one rating scale (e.g. Overall introduction with patient: *good, adequate or poor?*). Use of such rating scales can improve the reliability of an OSCE.¹¹

Cost efficiency of OSCEs

OSCEs are expensive and sophisticated forms of assessment. They are highly resource-dependent and require contributions from a large number of individuals. For example, a 16 station OSCE for over 250 medical students could require in excess of 128 examiner days. Of course there are also patients, faculty staff and other supporting personnel required for the assessment. Considerable effort is also required prior to the OSCE. In terms of planning the logistics of the exam also in development of the stations and training of assessors and patients. Costs regarding equipment, venue hire, catering and other sundry costs also need to be taken into account. Given the current economic imperative on academic institutions to make cost savings, there has never been a greater need to rationalise resources used in assessment. Later in this article I will discuss sequential OSCEs and their potential to reduce the number of examiners slots required - whilst maintaining the reliability of the assessment.

Acceptability of OSCEs

OSCEs need to be acceptable by all stakeholders. Therefore it is important to seek feedback from candidates, examiners and patients involved in the OSCE. Future employers of the candidates also need to have an active role. Given the perceived unfairness of the *long case*, OSCEs have become widely accepted and popular in undergraduate medical education.^{4, 8}

In OSCEs, all candidates should experience the same assessment experience and conditions. Inevitably there is potential for variation in OSCEs - for example between different circuits of the same OSCE and between different examiners.¹¹ The GMC have highlighted this issue and emphasise the importance of institutions paying special attention to assessor recruitment, training and monitoring.¹²

SETTING THE PASSING STANDARD IN OSCEs

To establish creditable standards, faculty must use a systematic approach in gathering expert judgments about acceptable levels of competency.¹⁵⁻¹⁶ To ensure the integrity and fairness of such passing scores, several standard setting procedures have been developed.¹⁶ Norm referenced (or *relative*) methods of standard setting are used when a fixed proportion of candidates are required to pass. In such methods of standard setting, competent candidates may fail to progress if the cohort are of above average ability. Therefore norm referencing methods of standard setting are generally unacceptable in undergraduate medical OSCEs. Methods that define a cut-off score, thereby identifying candidates who are competent and eligible for progression, are preferred in undergraduate OSCEs - i.e. criterion (or *absolute*) referencing. The borderline regression (BLR) method is a popular criterion-referenced method of setting a passing standard in OSCEs. The BLR method is generally considered robust and defensible.^{11, 14, 17-19}

In the BLR method - assessors directly observe candidates performing the clinical task in each station. They score the various components of the clinical task on the predefined

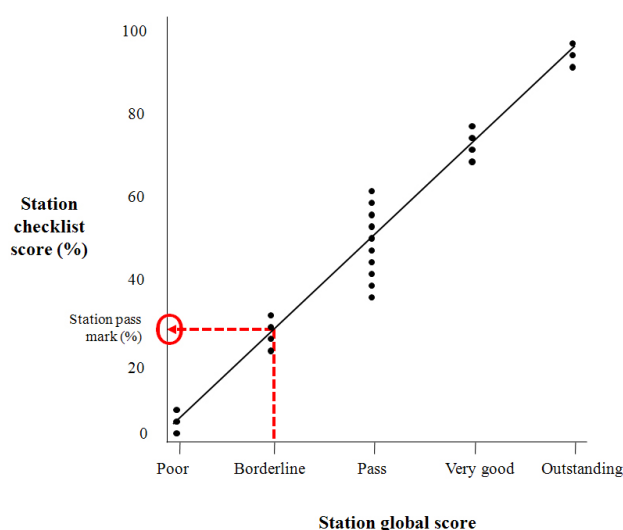


Fig 3. Graphical representation of the borderline regression method of calculating an OSCE station pass mark (i.e. linear regression of station checklist scores on to global scores)

checklist. Assessors then provide a separate overall rating or a *global score* of the candidate's performance (for example: *Outstanding, very good, pass, borderline or fail*). The pass mark for each OSCE station is then calculated by statistically regressing candidates' *checklist scores* on *global scores* for each station (Figure 3).

The overall pass mark of the OSCE is calculated by aggregating the pass marks for each of the separate OSCE stations. Upward adjustments maybe made by using the Standard Error of Measurement (SEM). Making such an adjustment reduces the probability of passing an incompetent candidate.²⁰ However there is also a chance of failing an only-just competent candidate. Protecting patients from incompetent doctors would support the argument for making such adjustments.

ASSESSORS IN OSCEs

Assessors play a vital role in delivering a robust and fair OSCE. Ultimately the decision to pass or fail a candidate in an OSCE does not fall on one assessor but on the entire panel of assessors. In the United States simulated patients often act as assessors in OSCEs.¹⁴ However in the United Kingdom and other parts of the world, clinicians tend to examine in OSCEs.

There is an imperative that institutions ensure assessors are competent to undertake their role.^{6, 12, 14} The GMC set out clear recommendations of the roles and responsibilities of assessors.¹ The Academy of Medical Educators also set out professional standards of *good educators* involved in assessment.²¹ Ideally the only variation in OSCEs should be due to candidates' performance and not due to any assessor effects or bias. Therefore in order for assessors to carry out their role consistently, they require training and feedback on their judgements and behaviour.¹² Most institutions now have established training programmes for OSCE assessors. At Queen's University Belfast we also supplement assessor training with an online learning module (www.med.qub.ac.uk/OSCE).²² This online training package outlines the roles and responsibilities of an OSCE assessor. Users are also provided the opportunity to practice scoring on an OSCE *checklist* and awarding *global scores* using online videos. In an anonymised fashion they can calibrate their decisions by comparing their awarded scores with that of their peers. However there remains a need for research in this area particularly on the effect that training has on assessor variance in OSCEs.²³

PATIENTS IN OSCEs

Most OSCE stations allow the observation of candidates interacting with patients. Patients may be either *real* or *simulated*. Real patients provide the opportunity to assess candidates' ability to examine for actual clinical features (e.g. auscultation for a cardiac murmur or examining a thyroid goitre). There are, however, significant issues regarding the use of real patients in OSCEs.²⁴ Firstly, OSCEs are demanding and have the potential to cause discomfort to a patient after being repeatedly examined by a large cohort of students (e.g. knee examination in a patient who has osteoarthritis). Furthermore real patients, and their clinical features, are often difficult to standardise - which can lead to candidates experiencing differences in OSCEs. Because of these challenges there ultimately has been a reduction in the use of real patients in undergraduate medical OSCEs.²⁵

Such a reduction in the use of real patients in OSCEs appears to influence some medical students' learning behaviours. In a recent survey of Final MB medical students, patients with *cardiac murmurs* and *pulmonary fibrosis* were predicted as the 'most likely' types of real patient cases that would occur in a graduating OSCE.²⁵ These predictions were based on the notion that such clinical cases were easy to standardise across different examination venues and amenable to repeated examinations. Such strategic predictions appear to influence students in their learning and encourage them to ignore 'less likely' cases in their clinical training. Faculty need to meet the challenges of using real patients in OSCEs and widen their participation.

Simulated patients (i.e. individuals without actual clinical features) are more commonly used in OSCEs. They can be used in different formats in order to portray a clinical scenario. For example they may be given a script of the symptoms of a patient who presents with acute coronary syndrome. Candidates then have to elicit the clinical history from the simulated patient. Scripts that are based on actual patients' accounts of their condition enhance the validity and patient centeredness of the OSCE station.²⁶ Simulated patients can also facilitate the assessment of candidates' physical examination skills (e.g. performing an abdominal examination). Simulated patients can also mimic certain clinical signs (e.g. a visual field defect or 'tenderness' in their right iliac fossa). However the potential range of signs that can adequately be reproduced are limited. Such clinical sign simulation requires effective training of simulated patients in order for them to portray the signs consistently.



Fig 4. Example of part-task simulation. A venepuncture manikin arm is attached to a simulated patient. Candidates in this station are asked to obtain a venous blood sample from the manikin arm but also interact and explain the procedure to the simulated patient.



Fig 5. Example of an inexpensive method of high fidelity simulation in an OSCE station. In this station - a temporary transfer tattoo of a malignant melanoma is placed on a simulated patient. Candidates are asked to interact with the patient, assess the 'skin lesion' and explain the potential diagnosis to the patient.

Increasingly the use of manikins and other technical equipment, in combination with simulated patients, are being used in OSCEs. For example attaching a venepuncture manikin arm to a simulated patient (Figure 4).

Such *hybrid* or *part-task* simulation not only allows for the assessment of the technical aspects of the clinical skill but also the humanistic dimensions of the encounter. Another example of such enhanced simulation include the use of high fidelity transfer tattoos of skin lesions.²⁷ The use of temporary tattoos can allow candidates to be assessed on their ability to diagnosis a skin lesion in a more realistic and patient-centred context (for example a high fidelity transfer tattoo of a malignant melanoma).

The Ventriloscope® is an electronic stethoscope that can realistically and consistently simulate 'abnormal' auscultatory findings.²⁸ Such technology appears to enhance validity within an OSCE setting.²⁹

CONTEMPORANEOUS ISSUES RELATING TO OSCEs

Patient ratings on candidates' performance in OSCEs.

Where appropriate, patients are often asked to rate a candidate's performance in an OSCE station.^{26,30} For example at Queen's University Belfast we pose our simulated patients with the following statement 'I would be happy to come back and discuss my concerns with this student again'. Simulated patients then provide a response using the following scale (*Strongly agree, agree, just agree, neutral or disagree*). Such ratings tend to focus on the humanistic aspects of the clinical encounter (e.g. attentiveness, empathy and rapport). There are a number of reasons why simulated patients are asked to rate candidates' performances. Not only does it highlight the importance of patient-centred care to our students, it also promotes simulated patients engagement in the assessment process. Furthermore, including simulated patients ratings to assessors *checklist* scores can potentially enhance the psychometric reliability of an OSCE.³¹ Simulated patients' ratings may also be used as a separate progression criteria for candidates in an OSCE (eg. regardless of the total OSCE

score, a candidate may fail to progress if a minimum number of simulated patients do not rate their performance as being satisfactory). Such a process requires effective training and quality control of simulated patients and their decisions.

Sequential OSCEs

As outlined previously in this paper, OSCEs are complex and expensive forms of assessment. In recent times sequential OSCEs have been developed so that reliability of the assessment is maintained but resources are targeted where they are needed the most i.e. the pass / fail divide.³² In a sequential OSCE, candidates go through an OSCE with a reduced number of stations (for example a 10 rather than a 16 station OSCE). The BLR is used to determine the cut score in this OSCE. However an upward adjustment of 2 or more SEMs are made to this pass mark. This invariably will produce a larger cohort of candidates who don't meet the standard. Within this group of candidates there are those that are truly *incompetent* and others who are truly *competent*. This group of candidates then go through an extended OSCE (e.g. a further 6 stations). Therefore the overall reliability of correctly identifying those students who are competent, in this small cohort of candidates, is maintained. In essence the OSCE does not have to be as reliable for all candidates, but focuses on those who are on the pass / fail boundary. Such an OSCE design requires fewer examiner days - which is of course more cost effective.

Quarantining ('corralling') in OSCEs

OSCEs often span the course of a day or more. With such practice there is potential for OSCE content to be leaked between different cohorts of candidates sitting the same examination. However there is a general consensus in the literature that such conduct does not have any significant statistical bearing on candidates' performance in OSCEs.³³⁻³⁴ OSCEs assess *showing* rather than *knowing* skills. Therefore the notion is that there is insufficient time to rehearse a skill in order to obtain any advantage.³⁴ Nonetheless such violations of OSCE content can potentially endanger the integrity and creditability of the assessment process. Therefore some institutions quarantine candidates between different sittings of the same OSCE (i.e. following an earlier sitting of an OSCE, candidates are placed in a holding area without access to their mobile phones or other electronic devices - until the next cohort of candidates have finished the OSCE).

Serious concern ('yellow card') reporting systems

A criticism of OSCEs is that a candidate can be incompetent in a particular skill, but can still pass the overall OSCE due to compensation from their performance in other stations. In response to this criticism, a number of institutions, including Queen's University Belfast, have developed a serious concern or 'yellow card' reporting system in their OSCEs. Such a system represents a qualitative mechanism of providing feedback to a candidate (and faculty) about their performance in an OSCE. Issues that would warrant a serious concern report include unprofessional practice (e.g. *being rough with a patient*) or unsafe actions that could potentially cause harm to a patient in clinical practice (e.g. *administration of an incorrect and dangerous drug*). In such significant situations candidates are asked to meet with faculty in order

to critically review the event. Before such candidates can progress on with their studies they are required to go through a remedial process until they have satisfactorily demonstrated competency in that particular skill. Future research is required to examine the predictive validity of serious concerns reports on future student performance in clinical practice.

CONCLUSION

Since their original development, OSCEs have become one of the main methods of assessing clinical competence in undergraduate medical education. Without question, OSCEs are more reliable than traditional methods of assessing clinical competence such as the *long case*. However they are not without their weaknesses. The high reliability of OSCEs is often at the expense of their validity. However with increased validity evidence, OSCEs have become more sophisticated and are portraying more realistic clinical scenarios. Used in combination with other methods of assessing clinical competency the shortcomings of OSCEs can be minimised. If correctly designed OSCEs can have a beneficial impact on medical students learning and future performance.

The author has declared no conflict of interest.

REFERENCES

1. General Medical Council. Good medical practice: duties and responsibilities of doctors. London: General Medical Council; 2009.
2. Harden RM, Gleeson FA. Assessment of clinical competence using an objective structured clinical examination (OSCE). *Med Educ.* 1979; **13**(1): 41-54.
3. Newble D. Techniques for measuring clinical competence: objective structured clinical examinations. *Med Educ.* 2004; **38**(2):199-203.
4. Ponnamperna GG, Karunathilake IM, McAleer S, Davis MH. The long case and its modifications: a literature review. *Med Educ.* 2009; **43**(10):936-41.
5. Epstein RM. Assessment in medical education. *N Engl J Med.* 2007; **356**(4): 387-96.
6. General Medical Council. Tomorrow's doctors: outcomes and standards for undergraduate medical education. 2nd ed. London: General Medical Council; 2009.
7. Miller GE. The assessment of clinical skills/competence/performance. *Acad Med.* 1990; **65**(9 Suppl): S63-7.
8. Wass V, Van der Vleuten C, Shatzer J, Jones R. Assessment of clinical competence. *Lancet.* 2001; **357**(9260): 945-9.
9. Van der Vleuten CP. The assessment of professional competence: developments, research and practical implications. *Adv Health Sci Educ.* 1996; **1**(1): 41-67.
10. Swanson DB. A measurement framework for performance based tests. In: Hart I, Harden R, editors. Further developments in assessing clinical competence. Montreal: Can-Heal; 1987. p. 13 - 45.
11. Pell G, Fuller R, Homer M, Roberts T. How to measure the quality of the OSCE: A review of metrics - AMEE guide no.49. *Med Teach.* 2010; **32**(10): 802-11.
12. General Medical Council. Assessment in undergraduate medical education: advice supplementary to Tomorrow's Doctors. London: General Medical Council; 2009. Available online from: http://www.gmc.uk.org/Assessment_in_undergraduate_web.pdf_38514111.pdf . [Last accessed August 2011].
13. Newble DI, Jaeger K. The effect of assessments and examinations on the learning of medical students. *Med Educ.* 1983; **17**(3): 165-71.
14. Boursicot K, Etheridge L, Setna Z, Sturrock A, Ker J, Smee S, et al.

- Performance in assessment: consensus statement and recommendations from the Ottawa conference. *Med Teach*. 2011; **33**(5): 370-83.
15. Livingstone SA, Zieky MJ. Passing Scores: a manual for setting standards of performance on educational and occupational tests. Princeton: Educational Testing Service; 1982.
 16. Norcini JJ. Setting standards on educational tests. *Med Educ*. 2003; **37**(5): 464-9.
 17. Wood TJ, Humphrey-Murto SM, Norman GR. Standard setting in a small scale OSCE: a comparison of the Modified Borderline-Group Method and the Borderline Regression Method. *Adv Health Sci Educ Theory Pract*. 2006; **11**(2): 115-22.
 18. Kramer WM, Muijtjens A, Jansen K, Dusman H, Tan L, van der Vleuten C. Comparison of a rational and an empirical standard setting procedure for an OSCE. Objective structured clinical examination. *Med Educ*. 2003; **37**(2): 132-9.
 19. Smee SM, Blackmore DE. Setting standards for an objective structured clinical examination: the Borderline Group Method gains ground on Angoff. *Med Educ* 2001; **35**(11): 1009-10.
 20. Streiner DL, Norman GR. Health measurement scales: a practical guide to their development and use. 3rd ed. Oxford: Oxford Medical publications; 2003.
 21. Academy of Medical Educators. Professional standards. London: Academy of Medical Educators; 2009. Available online from: <http://bit.ly/n0A5TG> [Last accessed August 2011]
 22. Centre for Medical Education. Queen's University Belfast [Internet]. OSCE examiner training and development. c2010-. Queen's University Belfast. Available online from: www.med.qub.ac.uk/OSCE. [Last accessed August 2011].
 23. Boursicot KA, Roberts TE, Pell G. Using borderline methods to compare passing standards for OSCEs at graduation across three medical schools. *Med Educ*. 2007; **41**(11): 1024-31.
 24. Collins JP, Harden RM. AMEE Medical Education guide No. 13: real patients, simulated patients and simulations in clinical examinations. *Med Teach*. 1998; **20**(6): 508-21.
 25. Gormley GJ, McCusker D, MA Booley, McNeice A. The use of real patients in OSCEs: Survey of medical students' opinions and predictions. *Med Teach* [Forthcoming 2011].
 26. Nestel D, Kneebone R. Perspective: authentic patient perspectives in simulations for procedural and surgical skills. *Acad Med*. 2010; **85**(5):889-93.
 27. Langley RGB, Tyler SA, Ornstein AE, Sutherland AE, Mosher LM. Temporary tattoos to simulate skin disease: report and validation of a novel teaching tool. *Acad Med*. 2009; **84**(7): 950-3.
 28. Lecat P. Lecat's VentriloScope. [Internet]. Ohio: North Eastern Ohio Universities Colleges of Medicine. Available online from: <http://ventriloSCOPE.com/>. [Last accessed July 2011].
 29. Verma A, Bhatt H, Booton P, Kneebone R. The VentriloScope® as an innovative tool for assessing clinical examination skills: Appraisal of a novel method of simulating auscultatory findings. *Med Teach*. 2011; **33**(7): e388-96.
 30. Kilminster S, Roberts T, Morris P. Incorporating patients' assessment into objective structured clinical examinations. *Educ Health (Abingdon)*. 2007; **20**(1): 6.
 31. Homer M, Pell G. The impact of the inclusion of simulated patient ratings on the reliability of OSCE assessments under the borderline regression method. *Med Teach*. 2009; **31**(5): 420-5.
 32. Cookson J, Crossley J, Fegan G, McKendree J, Mohsen A. A final clinical examination using a sequential design to improve cost effectiveness. *Med Educ* 2011; **45**(7):741-7.
 33. Colliver JA, Barrows HS, Vu NV, Verhulst SJ, Mast TA, Travis TA. Test security in examinations that use standardized-patient cases at one medical school. *Acad Med*. 1991; **66**(5):279-82.
 34. Swanson DB, Clauser BE, Case SM. Clinical skills assessment with standardized patients in high-stakes tests: a framework for thinking about score precision, equating, and security. *Adv Health Sci Educ Theory Pract*. 1999; **4**(1):67-106.